



<b>Title</b>	Combining biomarker and food intake data: calibration equations for citrus intake
<b>Authors(s)</b>	D'Angelo, Silvia, Gormley, Isobel Claire, McNulty, Breige A., Brennan, Lorraine, et al.
<b>Publication date</b>	2019-10
<b>Publication information</b>	D'Angelo, Silvia, Isobel Claire Gormley, Breige A. McNulty, Lorraine Brennan, and et al. "Combining Biomarker and Food Intake Data: Calibration Equations for Citrus Intake." Oxford University Press, October 2019. <a href="https://doi.org/10.1093/ajcn/nqz168">https://doi.org/10.1093/ajcn/nqz168</a> .
<b>Publisher</b>	Oxford University Press
<b>Item record/more information</b>	<a href="http://hdl.handle.net/10197/11209">http://hdl.handle.net/10197/11209</a>
<b>Publisher's statement</b>	This is a pre-copyedited, author-produced PDF of an article accepted for publication in The American Journal of Clinical Nutrition following peer review. The definitive publisher-authenticated version Silvia D'Angelo, Isobel Claire Gormley, Breige A McNulty, Anne P Nugent, Janette Walton, Albert Flynn, Lorraine Brennan, Combining biomarker and food intake data: calibration equations for citrus intake, The American Journal of Clinical Nutrition, Volume 110, Issue 4, October 2019, Pages 977–983, is available online at: <a href="https://doi.org/10.1093/ajcn/nqz168">https://doi.org/10.1093/ajcn/nqz168</a>
<b>Publisher's version (DOI)</b>	<a href="https://doi.org/10.1093/ajcn/nqz168">10.1093/ajcn/nqz168</a>

Downloaded 2026-05-01 23:38:05

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd\_oa)



© Some rights reserved. For more information

**COMBINING BIOMARKER AND FOOD INTAKE DATA: CALIBRATION EQUATIONS FOR CITRUS INTAKE**

Silvia D'Angelo<sup>1,2</sup>, Isobel Claire Gormley<sup>2</sup>, Breige A McNulty<sup>1</sup>, Anne P Nugent<sup>3</sup>, Janette Walton<sup>4,5</sup>, Albert Flynn<sup>4</sup>, and Lorraine Brennan<sup>1</sup>.

Names for PubMed indexing: D'Angelo, Gormley, McNulty, Nugent, Walton, Flynn, Brennan.

<sup>1</sup> Institute of Food and Health, School of Agriculture and Food Science, University College Dublin, Dublin, Ireland

<sup>2</sup> School of Mathematics and Statistics, Insight Centre for Data Analytics, University College Dublin, Dublin, Ireland

<sup>3</sup> Institute for Global Food Security, School of Biological Sciences, Queens University Belfast, Northern Ireland.

<sup>4</sup> School of Food and Nutritional Sciences, University College Cork, Cork, Ireland

<sup>5</sup> Dept. Biological Sciences, Cork Institute of Technology, Cork, Ireland

Address correspondence to L Brennan, Institute of Food and Health, Belfield, University College Dublin, Dublin 4, Ireland. E-mail: [lorraine.brennan@ucd.ie](mailto:lorraine.brennan@ucd.ie) ; Phone: 00 353 1 7162811

Supported by a research grant from The European Research Council ERC (647783) and from a Science Foundation Ireland grant (SFI/12/RC/2289).

Running heading: Calibration equations for citrus intake

Abbreviations: MSE (Mean Squared Error),  $R^2$  (R squared).

## **Abstract**

### **Background**

Measurement error associated with self-reported dietary intake is a well-documented issue. Combining biomarkers of food intake and dietary intake data is a high priority.

### **Objectives**

The objective was to develop calibration equations for food intake, illustrated with an application on citrus intake. Further, a simulation-based framework was developed to determine the portion of biomarker data needed for stable calibration equations estimation in large population studies.

### **Design**

Calibration equations were developed using mean daily self-reported citrus intake (4 days semi-weighed food diaries) and biomarker-derived intake (urinary proline betaine biomarker) data from participants (n=565) as part of a cross-sectional study. Different functional specifications and biomarker transformations were tested to derive the optimal calibration equations specification. The simulation study was developed using linear regression for the calibration equations. Stability in the calibration equations estimation was investigated for varying portions of biomarker and intake data “qualities”.

### **Results**

With citrus intake, linear regression on non-transformed biomarker data resulted in the optimal calibration equations specification and produced good quality predicted intakes. The lowest MSE (14354) corresponded to a linear regression model, defined with biomarker-derived estimates of intakes on the original scale. Using this model in a sub-population without biomarker data resulted in an average mean citrus intake of  $81 \pm 66$  g/day. The simulation study suggested that in large population studies, biomarker data on 20-30% of the

subjects are required to guarantee stable estimation of calibration equations. The paper is accompanied by a web application (“Bio-Intake”), developed to facilitate measurement error correction in self-reported mean daily citrus intake data.

### **Conclusion**

Calibration equations proved to be a useful instrument to correct measurement error in self-reported food intake data. The simulation study demonstrated that the use of food intake biomarkers may be feasible and beneficial in the context of large population studies.

### **Keywords**

Biomarkers, calibration equations, measurement error, citrus, proline betaine

## INTRODUCTION

Problems surrounding our ability to accurately measure dietary intake have been well documented in recent years in the literature (1-3). To overcome the measurement error associated with self-reported dietary intake the concept of dietary biomarkers has emerged. These biomarkers can deliver objective measures of intake and examples include urinary sodium and nitrogen for salt and protein intake respectively (1). Furthermore, there is evidence that use of biomarkers to develop calibrated equations to correct self-reported intake for estimates for protein and energy intake can positively impact on the ability to examine diet-disease associations: for example inclusion of biomarker calibrated data in the Women's Health Initiative cohorts allowed for disease associations to be revealed that otherwise would not have been identified (4, 5). Using urinary nitrogen and doubly labelled water as biomarkers to calibrate intakes of protein and energy in a group of postmenopausal women revealed associations between protein and energy intake with diabetes risk. It is noteworthy that these associations were not identified in the uncalibrated data.

The majority of the work to date in this field has focused on energy and protein intake using urinary nitrogen and doubly labelled water as recovery biomarkers. However, recent data has demonstrated that a series of biomarkers including carotenoids, tocopherols, folate, vitamin B12, and phospholipid fatty acids performed as well as established energy (doubly labelled water) and protein (urinary nitrogen) biomarkers in representing nutrient intake (6). This highlights the potential of widening the nutrients that could be corrected for measurement error using biomarkers and demonstrates that estimating intakes is not limited to recovery biomarkers. Furthermore, data emerging from the application of metabolomics to biomarker discovery has revealed that food intake biomarkers may be used to estimate food intake opening up the possibility of development of calibration equations for intake of foods (7-9).

The objective of the present work was to develop calibration equations for citrus intake using urinary proline betaine as the biomarker. Furthermore, we demonstrate through a simulation study the required portion of biomarker measures necessary to enable development of calibration equations in a large population group thus laying the foundations for future work. Finally, a web application (“Bio-Intake”) was developed to perform correction of measurement error in self-reported mean daily citrus intake data, based on the findings of this paper (<https://adiet.shinyapps.io/Bio-Intake>.)

## **METHODS**

### **STATISTICAL ANALYSIS**

#### **Development of calibration equations for estimation of citrus intake**

All data used for the calibration equations was obtained from previously reported studies (7). The selection process for the data is illustrated through a flow chart in the supplementary material (**Supplementary Figure 1**) and the key variables are self-reported citrus intake and urinary proline betaine. Briefly, dietary data was collected, over four consecutive days, using a four day semi-weighed food diary. Mean daily citrus intake (average citrus intake based on the four days of recording) was computed for the total citrus food group. The biomarker proline betaine has been previously presented a good biomarker of citrus intake with correlations of 0.9 with actual citrus intake. Furthermore, this has been demonstrated in different study settings by multiple research teams (9-12). Our previous work indicated that urinary proline betaine levels accurately estimated total citrus intake in a free living population group (7).

The calibration equations method relies on a central assumption, that both self-reported and biomarker-derived estimates of intake of study participant  $i$  are linearly related to the so-

called “true” intake  $X_i$ , with the relations being perturbed by Gaussian errors. However, as true intakes are unknown, directly estimating such relationships is not possible. The calibration method overcomes this issue by deriving a conditional prediction of  $X_i$ ,  $\hat{X}_i$ , given the biomarker-derived estimates of intake  $M_i$  (13). Such conditional predictions are derived from the calibrated intakes  $W_i$ , i.e. the self-reported intakes adjusted for measurement error:

$$M_i = \beta_0 + \beta_1 W_i + \epsilon_i ; \quad \hat{X}_i = \hat{\beta}_0 + \hat{\beta}_1 W_i$$

The data analyzed herein refer to a study conducted on  $N = 1500$  subjects, indexed by  $i = 1, \dots, N$  and collectively denoted  $A_N$ , where food consumption data was recorded using 4 day semi weighed food diaries. Mean daily citrus intake  $W_i$  (g/day) was computed for each subject. Also, biomarker-derived estimates of citrus intake  $M_i$  were computed for a subsample of  $A_N$ , here denoted  $A_n$ , with  $n = 565$ . Self-reported data and biomarker estimated intake data for these  $n$  subjects were used to develop calibration equations. The estimated calibration equations were then employed to correct self-reported mean daily citrus intake  $W_i$  for the remaining  $(N - n) = n^* = 935$  subjects, collectively denoted  $A_{n^*}$ .

The supplementary material provides a more detailed description of the calibration method. Common problems that may be encountered with food intake data include zero inflation (14), due the presence of non-consumers, and non-Gaussianity of the errors  $\epsilon_i$ . Potentially, such issues could result in poor correction of measurement error. In the literature, possible non-Gaussianity of the errors has been addressed using different transformations of the biomarker-derived estimates of intake,  $M_i$ , such as the log transformation (15) or Box-Cox transformations (16). In previous work (16), zero-inflated data have been commonly modeled with the assistance of covariates. However, when such subject-specific information is not available, such an approach cannot be pursued.

In the present work, transformations of the biomarker-derived intake data have been considered for the estimation of calibration equations, as well as original scale biomarker-derived estimates of intake, to address potential non-Gaussianity of subject-specific errors. Additionally, to address the issue of zero-inflated data in the absence of covariates, zero-inflated models have been explored and compared with a standard linear regression model, to select the optimal model specification for the calibration equations. The average mean squared error (MSE) between biomarker-derived estimates of intake and calibrated mean-daily self-reported intake has been used to compare different model specifications and different transformations of the biomarker-derived intake data. The supplement presents more detail on the different models explored and gives a thorough discussion of the results obtained.

#### **Estimation of the required portion of biomarker data in large studies**

Acquisition of biomarker derived estimates of intake for all subjects may not be feasible in large population studies, due to high costs and limited biological samples availability. In order to estimate the required quantity of biomarker data needed in large studies a simulation study has been conducted. A set of  $P = 100$  datasets of  $N = 100,000$  subjects was created. For each dataset  $p$  ( $p = 1, \dots, P$ ) biomarker-derived estimates of intake and self-reported intake data have been simulated, according to the results obtained from the analysis of mean-daily citrus intake data. Specifically, true mean-daily citrus intake values (g/day) have been simulated from a Gaussian distribution, with mean and variance fixed as the empirical mean and variance of the calibrated mean-daily self-reported citrus intakes. Self-reported and biomarker-derived intakes (g/day) have been then constructed from the simulated true intakes, following a classical calibration equations framework (17). To represent different situations that could be encountered with real data, simulated self-reported and biomarker

intakes may either have moderate or high variability around the “true intake”. Such variation is regulated by the parameters  $\alpha_W$  and  $\alpha_M$  for self-reported and biomarker-derived intakes respectively. Also, different strengths of the association between self-reported and true intakes have been accounted for, to represent different self-reported data qualities. This feature is regulated by the parameter  $\beta_W$ ; further detail on parameters  $\alpha_W, \alpha_M$  and  $\beta_W$  are given in the supplement. Biomarker derived intakes are assumed to have been measured only for  $n_g$  of the  $N$  subjects, with  $n_g = gN$  and  $g = (0.01, 0.02, \dots, 0.99, 1)$ . For each value of  $n_g$ , simulated intakes were used to estimate calibration equations, with the corresponding parameters being estimated via cross-validation (100 repetitions, excluding 20% of the observations at each round). Standard deviations of the estimated parameters,  $sd(\widehat{\beta}_0)$  and  $sd(\widehat{\beta}_1)$ , were computed and compared for different values of  $n_g$ , i.e. for different sizes of the sub-sample of subjects for which biomarker data are available. A stabilization of the parameter estimates indicates that, after a given  $n_g$ , further increasing the number of subjects for which biomarker-derived intakes are computed would not yield a substantial improvement in measurement error correction. More details regarding the construction of the simulation study are given in the supplement, while the results are presented in the next section.

## RESULTS

### Estimating citrus intake using calibration equations

Summary statistics for self-reported and biomarker-derived mean daily citrus intakes (g/day) are reported in

**Table 1.** All examined measures of citrus intake have similar interquartile ranges and similar means,  $\overline{W_{AN}} = 75.42$  g/day and  $\overline{M_{An}} = 83.26$  g/day. Biomarker-derived estimates of

intake range in a wider interval ( $\text{Max}(M_{A_n}) = 1226.79$ ) compared to self-reported intakes ( $\text{Max}(W_{A_n}) = 845.75$ ). Non-consumers represent around 40% of subjects in self-reported intake data, while zero value biomarker-derived intakes are only 17% of the total.

Concordance between self-reported and biomarker-derived intakes was confirmed by the statistics summarizing the conditional distributions  $M_{A_n|W_{A_n}=0}$  and  $W_{A_n|M_{A_n}=0}$  in

Table 1. These distributions represent, respectively, biomarker-derived intake values corresponding to null self-reported intake values and vice-versa. Such measures of intake display good agreement, as both take low values when the other is null.

Summaries of the predicted intakes obtained using different model specifications and biomarker transformations are compared in

**Table 2.** Models using a log-transform of the biomarker-derived estimates of intake in the estimation of the calibration equations, or of a Box-Cox transformation with small  $\lambda$  parameter, result in high-variability in the calibrated intakes. Indeed, regardless of the functional specification of the calibration equations, in these cases the standard deviation of the calibrated intakes is always greater than 792, much larger than the standard deviations of the observed self-reported intakes ( $sd(W_{A_n^*}) = 107$ , see Table 1). Also, when using these biomarker transformations, poor agreement is found between the interquartile ranges of the observed self-reported intakes and the calibrated intakes. The calibrated intakes from the remaining model specifications and biomarker transformations range in similar intervals, closer to that of the self-reported mean daily citrus intakes. Among these, the calibrated intakes computed using a linear regression model for the calibration equations with original-scale biomarker-derived estimates of intake ( $M_i$ ) are those whose mean is closer to the mean of  $W_{A_n^*}$ . Indeed, a linear regression model was optimal, in terms of MSE, to estimate calibration equations of citrus intake ( $R^2 = 0.22$ ).

1 **Table 3** details the average MSE values associated with different model specifications and  
 2 transformations of biomarker-derived estimates of intake. Mean squared errors have been  
 3 computed between biomarker-derived estimates of intakes and calibrated mean-daily self-  
 4 reported intakes, computed with different specifications of the calibration equations. The  
 5 lowest MSE (14354) corresponds to a linear regression model, defined with biomarker-  
 6 derived estimates of intakes on the original scale,  $M_i$ . The preference for a simple model is  
 7 confirmed by the second lowest MSE, associated with a linear regression model, with unitary  
 8 parameter Box-Cox transformation of  $M_i$ , coinciding with a mere translation of the data.  
 9 Given the selected optimal linear model, mean-daily calibrated citrus intakes were computed  
 10 using the estimated calibration equations:

$$11 \quad \hat{X}_i = 33.60 + 0.63W_i$$

12 It is worth noting that the estimated intercept is close to the average biomarker-derived  
 13 intake for non-consumers (see  
 14 Table 1). **Figure 1** offers a comparison between calibrated and self-reported mean-daily citrus  
 15 intakes for the  $n^*$  subjects in the validation study. Overall, there is good agreement between  
 16 self-reported and calibrated mean-daily intakes, however, differences emerge at high intakes.  
 17 Self-reported non-consumers were assigned to a low baseline citrus consumption level,  
 18 corresponding to  $\hat{\beta}_0 = 33.60$ .

### 19 **What portion of biomarker data is required in large population studies?**

20 Sixteen distinct scenarios of simulated intake data were generated by considering different  
 21 combinations of variability of self-reported data and biomarker data around the true intake,  
 22 quantified through  $\alpha_W, \alpha_M = (0.5, 2)$ , and self-reported data quality levels, quantified  
 23 through  $\beta_W = (0.1, 0.5, 0.8, 1)$ . Among these, the group of “best-case” scenarios corresponds

24 to those with a  $\alpha_W = \alpha_M = 0.5$ , as both biomarker-derived intakes and self-reported intakes  
25 have very little variability around true intakes.

26

27 **Figure 2** reports the standard deviations of the estimated calibration equations slope  
28 parameter  $\hat{\beta}_1$  under such scenarios, for the different self-reported data quality levels.  
29 Standard deviations stabilize quickly, for relatively small values of  $n_g$  (the number of subjects  
30 for which biomarker-derived estimates of intake are available), regardless of the  $\beta_W$  value.  
31 However,  $\beta_W$  does influence the level at which the standard deviation stabilizes, with poor-  
32 quality self-reported data ( $\beta_W = 0.1$ ) leading to slightly higher variability in the estimate of  
33  $\beta_1$ . Similar results are found also for the set of “worst-case” scenarios (**Figure 3**), represented  
34 by the  $\alpha_W = \alpha_M = 2$  configuration, that is, when large measurement error is present in both  
35 measures of intake. Standard deviations tend to stabilize quite rapidly around  $n_g \approx 25000$ ,  
36 as happened in the “best-case” scenarios. However, this time the quality of self-reported data  
37 does not have a large impact. Indeed, there is little visible differences between the cases  
38  $\beta_W = 0.1$  and  $\beta_W = 1$ . The remainder of the sixteen scenarios represent less extreme, hybrid  
39 situations in terms of data quality and measurement error, with corresponding results in-line  
40 with those described for the two extreme scenarios. Similar conclusions were drawn for the  
41 estimated calibration equations intercept parameter  $\beta_0$ ; standard deviations for the  
42 estimates of the intercept parameter stabilized jointly with those of  $\beta_1$ , around the same  $n_g$   
43 values. More detail regarding the simulation study, including results and MSEs associated with  
44 the different simulation scenarios, are provided in the supplement. In general, the standard  
45 deviations of the calibration equation parameters tend to stabilize when biomarker data is  
46 obtained from 25% of the study population, in this case for 25000 out of the 100000  
47 subjects.

## 48 **DISCUSSION**

49 The present study presents an important advancement in the use of intake biomarkers in  
50 nutrition research. Using a well-established biomarker of citrus intake calibration equations  
51 were developed to enable correction of self-reported intake for measurement error.  
52 Development and demonstration of such an approach for food intake clearly shows how one  
53 could use food intake biomarkers in nutrition research. To date the literature is lacking real  
54 examples of the utility of such food intake biomarkers with the majority of studies stopping  
55 at an association between food intake and biomarker levels. Within this context, the use of  
56 urinary proline betaine to correct self-reported citrus intake is a clear example of the potential  
57 of the field of food intake biomarkers. Furthermore, we developed a framework for assessing  
58 the amount of biomarker data needed in a large population study to correct for measurement  
59 error in food intake. Interestingly, in the present simulation study biomarker data was  
60 necessary for around 20-30% of the population highlighting an important aspect in the  
61 translation of this work into large epidemiological studies.

62 To date, calibration equations for correction of measurement error have focused on nutrients  
63 and to the best of our knowledge there is no example in relation to specific food intake. As  
64 dietary guidelines move towards food-based guidelines it is imperative that we also have  
65 good assessment of food intake. Use of food intake biomarkers has gained increased traction  
66 in recent years with the promise of improvement in dietary assessment. The current  
67 development of calibration equations for citrus intake clearly paves the path forward in this  
68 regard. The ability to correct for measurement error for food intake has the potential to  
69 enhance our understating of the relation between food and disease. Previous literature  
70 demonstrated that correction of measurement error in nutrient intake using biomarkers led  
71 to the identification of relationships that were not observed without the calibration. Using

72 biomarkers to correct for measurement error, energy intake was associated with increased  
73 risk of breast cancer, all cancer and type 2 diabetes and protein intake was associated with  
74 increased risk of type 2 diabetes (4, 18-20). Building upon the present work, the next step will  
75 be to examine relationships between calibrated food intake and certain health parameters.  
76 An important aspect of our approach was the development of calibrated intake distributions  
77 which incorporate the uncertainty inherent in statistical modelling of measurement error.  
78 Such an approach results in calibrated intake distributions instead of single point-estimates.  
79 For researchers the potential distribution of intake delivers valuable information considering  
80 the uncertainty surrounding self-reported data. To aid the end user to implement the citrus  
81 intake calibration equations herein we developed a freely available web application, "Bio-  
82 Intake". Bio-intake is an open-source, easy to use web application designed to allow  
83 reproducibility of the methodology proposed in this paper.

84 Typically, for biomarker-based studies the biomarker data is acquired on a sub-population to  
85 develop the calibration equations. In this context we developed a framework for the  
86 assessment of the portion of food intake biomarker data necessary to develop accurate  
87 calibration equations. The framework advocates for examination of the quality of the self-  
88 reported data and the biomarker derived data. The framework was illustrated through a  
89 stimulation study where the portion of biomarker measures ( $n_g$ ) necessary to enable accurate  
90 development of calibration equations in a large population group ( $N$  individuals) was  
91 estimated. For this purpose, standard deviations of the calibration equations parameter  
92 estimates are analyzed. A stabilization of such standard deviations occurring for low values of  
93  $n_g$ ,  $n_g \ll N$ , would indicate that the collection of biomarker data for a small portion of the  
94 population of interest and that for the whole population are practically equivalent in terms  
95 of stability of calibration equation estimation.

96 Indeed, simulation results showed rapidly decreasing standard deviations, independently of  
97 the scenario considered, that is of the “observed data-quality”. Such promising results suggest  
98 that substantial improvements in the correction of self-reported intake data may be achieved  
99 even with little biomarker data. Indeed, collecting biomarker-derived estimates of intake for  
100 a larger portion of the population of interest may not lead to more precise calibration  
101 equations estimates. In the current example of citrus intake, measurement of the biomarker  
102 in 20-30% of the population was deemed sufficient. The simulation study was replicated in a  
103 small cohort setting ( $N = 5000$ ), leaving unchanged the results. Obtaining biomarker data  
104 from just a small portion of the population suggests that calibration equations can prove to  
105 be powerful tools in large cohort studies, where cost and limited biological samples are  
106 barriers to their development. In the literature, previous studies have used biomarker data  
107 in 1-3 % of the total larger cohort to assess relationships between intake and disease risk  
108 diabetes (4, 21).

109 The findings of the current simulation study were partly in line with previous large population  
110 studies for nutrients, in that it is not necessary to collect biomarker data for the whole  
111 population of interest to develop calibration equations. Nonetheless, simulation results  
112 showed a substantial decrease in the standard deviations of the calibration equations  
113 parameter estimates for portion of biomarker data in between  $n_g = 0.01N$  and  $n_g = 0.25N$ ,  
114 suggesting that a value of 1% may be suboptimal. Differences between subsequent standard  
115 deviations values can be used as proxies of improvement in the stability of calibration  
116 equations. Computing these differences in standard deviations for subsequent  $n_g$  values  
117 shows that they become (approximately) null for a portion of the population between 7% and  
118 18%, depending on the scenario considered. While relaxing the recommendation to 7% of  
119 cohort biomarker data allows for greater data-collection feasibility, it induces the risk of larger

120 uncertainty in calibration equation estimates. Moreover, the 7% of cohort biomarker data  
121 figure is based on the strong assumption of very good quality self-reported data; small  
122 deviations from this assumption can induce even larger uncertainty in the estimates.  
123 Although being potentially onerous to achieve, the more conservative recommendation of  
124 25% of cohort biomarker data is able to provide stable estimates independently of the self-  
125 reported data quality. Also, the simulation study was developed to analyze calibration  
126 equations estimates for food biomarkers and its results may not have direct applicability in  
127 the context of nutrient biomarkers.

128 Lastly, mean squared error values from the simulation study demonstrated that the quality  
129 of analyzed self-reported data is quite relevant, even with high quality biomarker data.  
130 Indeed, mean squared error values associated with “low data quality” scenarios ( $\beta_W = 0.1$ )  
131 were always much larger than those related to “higher data quality” ones ( $\beta_W = (0.5, 0.8, 1)$ ).  
132 More importantly, mean squared error distributions in “low data quality” scenarios did not  
133 vary particularly with different variability levels in biomarker data ( $\alpha_W = (0.5, 2)$ ). In  
134 conclusion, the simulation study findings suggest that calibration equations may be stably  
135 estimated even with a relatively small portion of biomarker data, but that the quality of  
136 calibrated intake predictions is strongly dependent on that of self-reported intake data. Large  
137 biomarker data samples and excellent quality biomarkers will not be enough to get reliable  
138 estimates of intake (calibrated intake) if the observed self-reported data are not reliable.

139 This study has a number of noteworthy strengths. The development of calibration equations  
140 was based on reliable biomarker data and good quality food intake data (estimate from 4-day  
141 semi weighted food diaries). Furthermore, the calibration equations framework developed  
142 (and the corresponding “Bio-Intake” web app) allows estimation of calibrated intake  
143 distributions, instead of single point estimates. Hence, the proposed framework allows to

144 better account for uncertainty inherent to calibrated intakes estimation: self-reported intake  
145 values are associated to a range of possible calibrated values, providing richer information on  
146 the quantity of interest.

147 A number of limitations need to be acknowledged. The proposed framework only accounts  
148 for self-reported intake data in the estimation of calibration equations. When available, the  
149 introduction of potentially pertinent subject-specific covariates in the modelling framework  
150 should be considered, as it may improve estimation of the calibrated intakes. However, as a  
151 drawback, controlling for different factors with covariates would lead to the estimation of less  
152 generalizable calibration equations, which would have narrower applicability. Moreover,  
153 future work may consider the development of citrus calibration equations in other  
154 populations, to test the generalizability of the results obtained in the present analysis. Finally,  
155 the simulation study provided results concerning the stability of calibration equations  
156 parameter estimates, in terms of portion of biomarker data needed for the estimation. This  
157 portion does not guarantee convergence of such estimates, that is, reliable estimated  
158 calibrated intakes. The results demonstrate that reliability of the calibrated intakes depends  
159 on the quality of both the biomarker-derived and self-reported intake data, which does not  
160 depend on the modelling framework adopted for the development of calibration equations.  
161 Good self-reported data quality should always be assessed separately and before performing  
162 estimating calibration equations.

163 In conclusion, the present study demonstrates the potential of food intake biomarkers in  
164 nutritional epidemiology. Future work will include development of calibration equations for  
165 other foods which will in turn improve the accuracy of dietary assessment. A further  
166 important aspect to be explored is whether the calibrated food intakes could be used in

167 disease-risk studies in larger cohorts in a similar fashion to the work performed with  
168 calibrated nutrient intake.

169

170

171 **Acknowledgements**

172 The authors' responsibilities were as follows—SD'A: conducted research, analyzed data and  
173 prepared manuscript draft; ICG: conducted research, analyzed data and contributed to  
174 manuscript correcting; BAM, APN, JW and AF: provided essential data; LB: designed  
175 research, conducted research, analyzed data and contributed to manuscript correcting. All  
176 authors read and approved the final manuscript. None of the authors had a conflict of interest.

## References

1. Bingham SA. Biomarkers in nutritional epidemiology. *Public Health Nutr* 2002;5(6a):821-7. doi: 10.1079/PHN2002368.
2. Kipnis V, Midthune D, Freedman L, Bingham S, Day NE, Riboli E, Ferrari P, Carroll RJ. Bias in dietary-report instruments and its implications for nutritional epidemiology. *Public Health Nutr* 2002;5(6a):915-23. doi: 10.1079/PHN2002383.
3. Dhurandhar NV, Schoeller D, Brown AW, Heymsfield SB, Thomas D, Sørensen TIA, Speakman JR, Jeansonne M, Allison DB. Energy balance measurement: when something is not better than nothing. *Int J Obes* 2015;39(7):1109-13. doi: 10.1038/ijo.2014.199.
4. Tinker LF, Sarto GE, Howard BV, Huang Y, Neuhouser ML, Mossavar-Rahmani Y, Beasley JM, Margolis KL, Eaton CB, Phillips LS, et al. Biomarker-calibrated dietary energy and protein intake associations with diabetes risk among postmenopausal women from the Women's Health Initiative. *Am J Clin Nutr* 2011;94(6):1600-6. doi: 10.3945/ajcn.111.018648.
5. Prentice RL, Huang Y, Kuller LH, Tinker LF, Horn LV, Stefanick ML, Sarto G, Ockene J, Johnson KC. Biomarker-calibrated energy and protein consumption and cardiovascular disease risk among postmenopausal women. *Epidemiology* 2011;22(2):170-9. doi: 10.1097/EDE.0b013e31820839bc.
6. Lampe JW, Huang Y, Neuhouser ML, Tinker LF, Song X, Schoeller DA, Kim S, Raftery D, Di C, Zheng C, et al. Dietary biomarker evaluation in a controlled feeding study in women from the Women's Health Initiative cohort. *Am J Clin Nutr* 2017;105(2):466-75. doi: 10.3945/ajcn.116.144840.
7. Gibbons H, Michielsen CJR, Rundle M, McNulty BA, Nugent AP, Gibney MJ, Brennan L, et al. Demonstration of the utility of biomarkers for dietary intake assessment; proline betaine as an example. *Mol Nutr Food Res* 2017;61(10).
8. Garcia-Perez I, Posma JM, Chambers ES, Nicholson JK, C. Mathers J, Beckmann M, Draper J, Holmes E, Frost G. An analytical pipeline for quantitative characterization of dietary intake: application to assess grape intake. *J Agric Food Chem* 2016;64(11):2423-31. doi: 10.1021/acs.jafc.5b05878.
9. Garcia-Perez I, Posma JM, Gibson R, Chambers ES, Hansen TH, Vestergaard H, Hansen T, Beckmann M, Pedersen O, Elliott P, et al. Objective assessment of dietary patterns by use of metabolic phenotyping: a randomised, controlled, crossover trial. *Lancet Diabetes Endocrinol* 2017;5(3):184-95. doi: 10.1016/S2213-8587(16)30419-3.
10. Pujos-Guillot E, Hubert J, Martin J-F, Lyan B, Quintana M, Claude S, Chabanas B, Rothwell JA, Bennetau-Pelissero C, Scalbert A, et al. Mass spectrometry-based metabolomics for the discovery of biomarkers of fruit and vegetable intake: citrus fruit as a case study. *J Proteome Res* 2013;12(4):1645-59. doi: 10.1021/pr300997c.
11. Heinzmann SS, Brown IJ, Chan Q, Bictash M, Dumas M-E, Kochhar S, Stamler J, Holmes E, Elliott P, Nicholson JK. Metabolic profiling strategy for discovery of nutritional biomarkers: proline betaine as a marker of citrus consumption. *Am J Clin Nutr* 2010;92(2):436-43. doi: 10.3945/ajcn.2010.29672.
12. Lloyd AJ, Beckmann M, Favé G, Mathers JC, Draper J. Proline betaine and its biotransformation products in fasting urine samples are potential biomarkers of habitual citrus fruit consumption. *Br J Nutr* 2011;106(6):812-24. doi: 10.1017/S0007114511001164.
13. Gormley IC, Bai Y, Brennan L. Combining biomarker and self-report dietary intake data: a review of the state of the art and an overview of concepts. *Stat Methods Med Res* (in press) 2019.
14. Zhang S, Midthune D, Guenther PM, Krebs-Smith SM, Kipnis V, Dodd KW, Buckman DW, Tooze JA, Freedman L, Carroll RJ. A new multivariate measurement error model with zero-inflated dietary data, and its application to dietary assessment. *Annals Appl Stat* 2011;5(2B):1456-87. doi: 10.1214/10-AOAS446.

15. Prentice RL, Mossavar-Rahmani Y, Huang Y, Van Horn L, Beresford SAA, Caan B, Tinker L, Schoeller D, Bingham S, Eaton CB, et al. Evaluation and comparison of food records, recalls, and frequencies for energy and protein assessment by using recovery biomarkers. *Am J Epidemiol* 2011;174(5):591-603. doi: 10.1093/aje/kwr140.
16. Kipnis V, Midthune D, Buckman DW, Dodd KW, Guenther PM, Krebs-Smith SM, Subar AF, Toozé JA, Carroll RJ, Freedman LS. Modeling data with excess zeros and measurement error: application to evaluating relationships between episodically consumed foods and health outcomes. *Biometrics* 2009;65(4):1003-10. doi: 10.1111/j.1541-0420.2009.01223.x.
17. Keogh RH, White IR. A toolkit for measurement error correction, with a focus on nutritional epidemiology. *Stat Med* 2014;33(12):2137-55. doi: 10.1002/sim.6095.
18. Prentice RL, Pettinger M, Tinker LF, Huang Y, Thomson CA, Johnson KC, Beasley J, Anderson G, Shikany JM, Chlebowski RT, et al. Regression calibration in nutritional epidemiology: example of fat density and total energy in relationship to postmenopausal breast cancer. *Am J Epidemiol* 2013;178(11):1663-72. doi: 10.1093/aje/kwt198.
19. Prentice RL, Shaw PA, Bingham SA, Beresford SAA, Caan B, Neuhauser ML, Patterson RE, Stefanick ML, Satterfield S, Thomson CA, et al. Biomarker-calibrated energy and protein consumption and increased cancer risk among postmenopausal women. *Am J Epidemiol* 2009;169(8):977-89. doi: 10.1093/aje/kwp008.
20. Zheng C, Beresford SA, Van Horn L, Tinker LF, Thomson CA, Neuhauser ML, Di C, Manson JE, Mossavar-Rahmani Y, Seguin R, et al. Simultaneous association of total energy consumption and activity-related energy expenditure with risks of cardiovascular disease, cancer, and diabetes among postmenopausal women. *Am J Epidemiol* 2014;180(5):526-35. doi: 10.1093/aje/kwu152.
21. Neuhauser ML, Tinker L, Shaw PA, Schoeller D, Bingham SA, Horn LV, Beresford SAA, Caan B, Thomson C, Satterfield S, et al. Use of recovery biomarkers to calibrate nutrient consumption self-reports in the Women's Health Initiative. *Am J Epidemiol* 2008;167(10):1247-59. doi: 10.1093/aje/kwn026.

Table 1

<b>CITRUS INTAKE</b>	$W_{A_N}$	$W_{A_n^*}$	$W_{A_n}$	$M_{A_n}$	$M_{A_n}   W_{A_n}=0$	$W_{A_n}   M_{A_n}=0$
<b>Mean</b>	75.42	72.41	80.39	83.26	30.52	24.18
<b>Standard deviation</b>	106.11	107.19	104.21	135.45	56.56	54.01
<b>Min</b>	0.00	0.00	0.00	0.00	0.00	0.00
<b>First quartile</b>	0.00	0.00	0.00	7.83	0.00	0.00
<b>Median</b>	31.88	22.00	43.75	34.48	12.81	0.00
<b>Third quartile</b>	117.31	113.00	122.75	98.86	35.07	7.12
<b>Max</b>	845.75	845.75	712.50	1226.79	441.20	302.50
<b>Zero intakes</b>	646 (43%)	424 (45%)	215 (38%)	94 (17%)	68 (32%)	68 (72%)

Summary statistics for the citrus mean daily intake data (g/day). Self-reported mean daily intakes are reported for the entire sample ( $W_{A_N}$ ), subjects for which no biomarker data is available ( $W_{A_n^*}$ ) and subjects for which biomarker-derived estimates of intake are computed ( $W_{A_n}$ ). The fourth column reports summary statistics for the biomarker-derived estimates of intake ( $M_{A_n}$ ).

Table 2

	MODEL	Linear regression model							Tobit model		Two-part model	
	BIOMARKER TRANSFORMATION	$M_i$	$\log(M_i)$	$BC(M_i, 0.1)$	$BC(M_i, 0.5)$	$BC(M_i, 1)$	$BC(M_i, 1.1)$	$BC(M_i, 1.5)$	$M_i$	$\log(M_i)$	$M_i$	$\log(M_i)$
CALIBRATED MEAN DAILY CITRUS INTAKE (g/day)	Mean	81	103	116	58	81	86	101	81	66	64	128
	Standard deviation	66	792	1299	69	66	69	82	67	180	77	1332
	Min	34	19	3	22	34	35	33	33	17	8	15
	First quartile	34	19	3	22	34	35	33	33	17	8	15
	Median	54	25	4	31	53	58	74	52	25	32	20
	Third quartile	107	52	19	66	107	115	150	108	62	94	45
	Max	563	24124	38554	891	563	549	527	553	4313	625	42066

Summary statistics for calibrated intake values using different calibration equations specifications, corresponding to various biomarker-derived intake data ( $M_i$ ) transformations and functional specifications. The calibrated intakes are those computed for the group of  $n^*$  subjects for which only self-reported mean daily citrus intake data (g/day) were available.

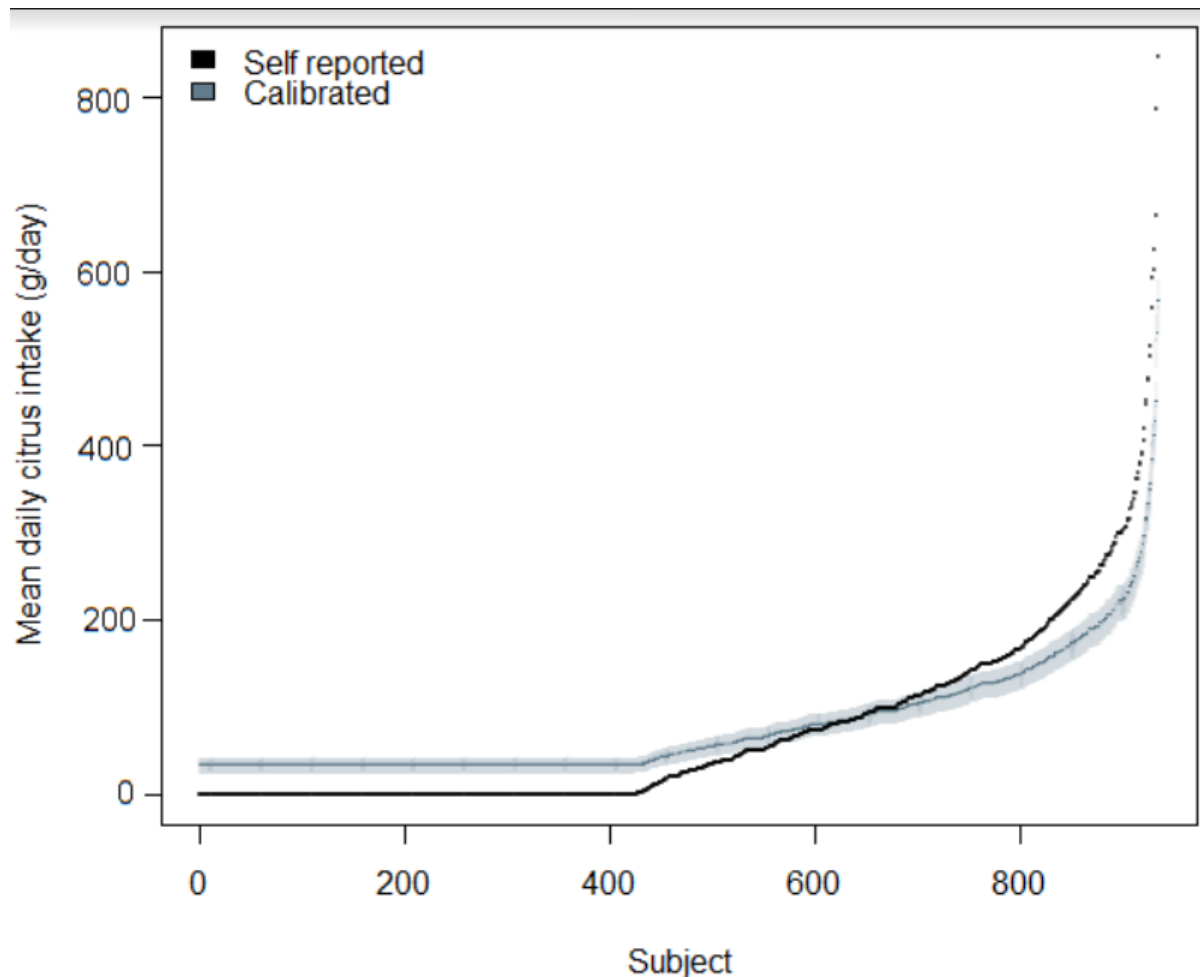
Table 3

<b>Average Mean Square Error</b>			
	<b>BIOMARKER TRANSFORMATION</b>		
<b>MODEL</b>	<b><math>M_i</math></b>	<b><math>\log(M_i + 1)</math></b>	<b><math>BC(M_i, 1)</math></b>
<b>Linear regression model</b>	14354	49078	14480
<b>Tobit model</b>	15034	89645	-
<b>Two-part model</b>	21504	21010	-

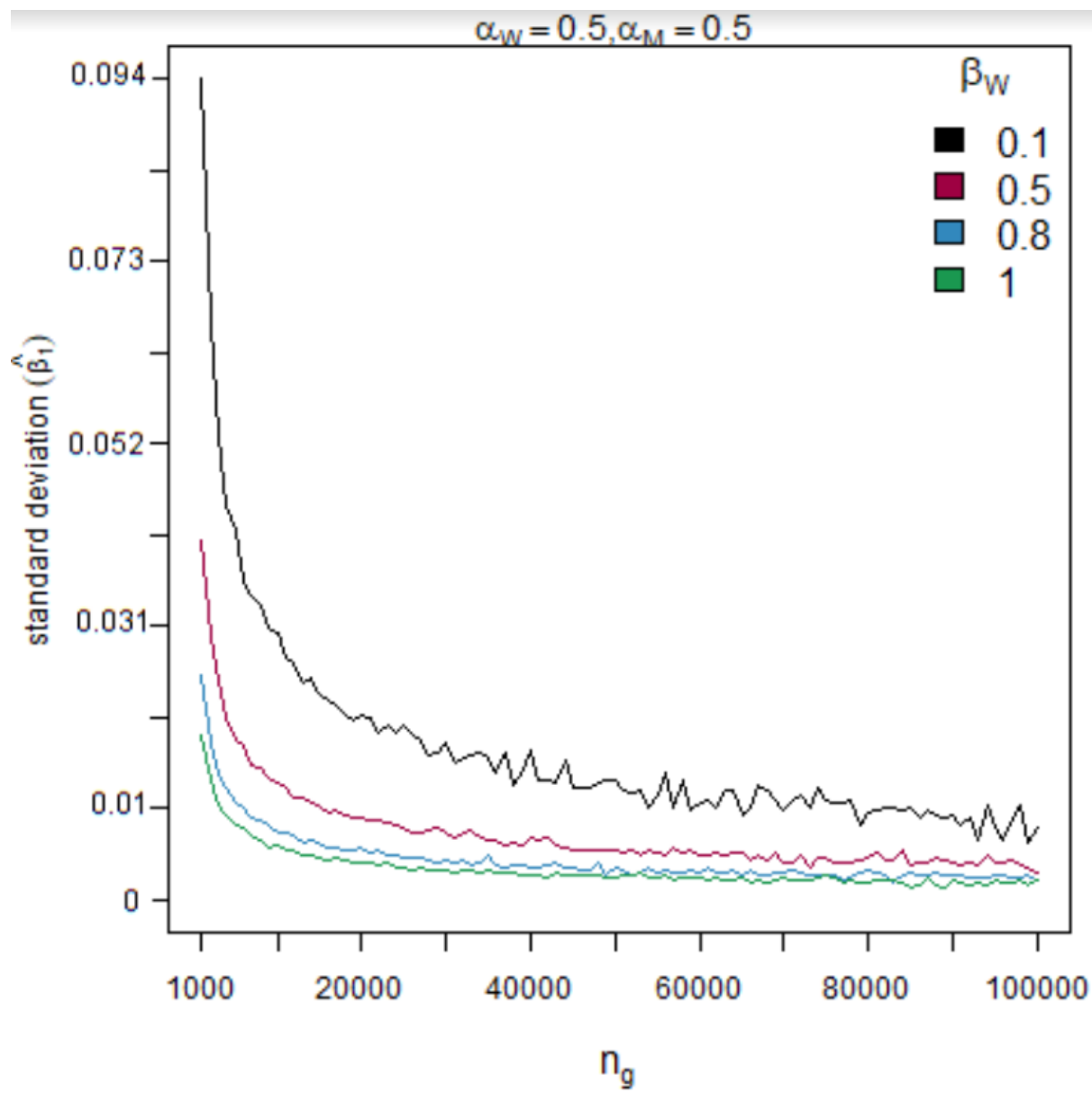
Average mean squared errors computed between calibrated intakes and biomarker-derived estimates of intakes ( $M_i$ ), for the  $n$  subjects in the study and different specifications of the calibration equations. The average MSE values for different specification of the Box-Cox  $\lambda$  parameter have been omitted, being all greater than the one obtained with  $\lambda = 1$ .

### Figure Legends

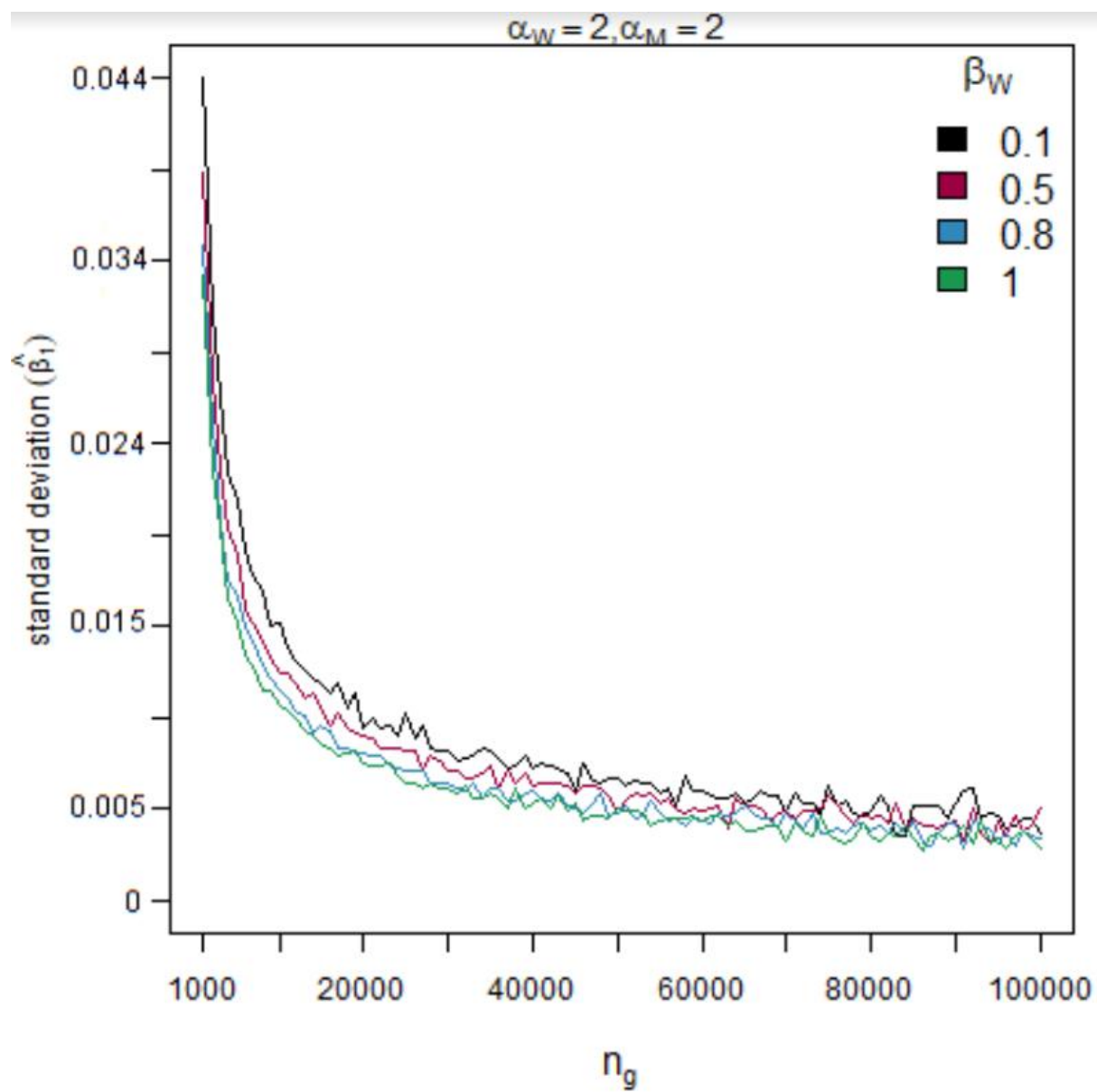
**Figure 1** Self-reported mean daily citrus intakes (g/day) and calibrated mean daily citrus intakes (g/day). For the calibrated intakes, 95% confidence intervals are reported (light grey). Intakes are ordered from lowest to highest values.



**Figure 2** Standard deviations of the estimated slope coefficient  $\beta_1$  are reported, for different sample sizes  $n_g$  and different values of the  $\beta_W$  parameter. These simulated scenarios represent the “best possible” ones in terms of data quality, as both self-reported and biomarker-derived estimates of intakes have low variability around the true intake values. Notice that the standard deviations tend to stabilize quickly, even in the “worst-case” scenario, that is when  $\beta_W = 0.1$ .



**Figure 3** Standard deviations of the estimated slope coefficient  $\beta_1$  are reported, for different sample sizes  $n_g$  and different values of the  $\beta_W$  parameter. These simulated scenarios represent the “worst possible” outcomes in terms of data quality, as both self-reported and biomarker-derived estimates of intakes have high variability around the true intake values. Notice that the standard deviations tend to stabilize quickly, even in the “worst-case” scenario, that is when  $\beta_W = 0.1$ .



## SUPPLEMENTARY MATERIAL

### COMBINING BIOMARKER AND FOOD INTAKE DATA: CALIBRATION EQUATIONS FOR CITRUS INTAKE

Silvia D'Angelo, Isobel Claire Gormley, Breige A McNulty, Anne P Nugent, Janette Walton, Albert Flynn, and Lorraine Brennan

#### THE CALIBRATION METHOD

The selection process for the data used to develop calibration equations is illustrated by **Supplementary Figure 4**.

##### How it works

The calibration method is a measurement error correction technique widely used in nutrition literature (1-3). The main assumption behind it states that self-reported intake data  $W_i$  are linearly related to true intakes  $X_i$ :

$$W_i = \mu_w + \beta_w X_i + \epsilon_{w,i} \quad (1)$$

for subjects  $i = 1, \dots, N$ . The relation between the observed data  $W_i$  and the latent quantity of interest  $X_i$  is mainly influenced by two factors: the slope coefficient  $\beta_w$  and the subject-specific errors  $\epsilon_{w,i}$ . The  $\beta_w$  parameter may be seen as a measure of the “quality” or the intensity of the relationship between self-reported and true intake data. Also, it expresses the direction of the relation between the two variables, with negative values of  $\beta_w$  indicating disagreement. Subject-specific errors  $\epsilon_{w,i}$  represent individual variability in self-reported intakes. Such errors are assumed to be normally distributed with zero mean and homoscedastic, i.e. with equal variances  $\sigma_w^2$ . Thus, the calibration method assumes there is no interaction between the subjects when reporting the intakes, a reasonable assumption in many nutrition studies, as reports are collected separately. Also, by assuming equal error variances  $\sigma_w^2$ , the calibration method implies that subjects are equally likely to misreport their consumption, hence all subjects are equally reliable.

The linear relation between self-reported and true intake data cannot be estimated, as the latter is not observed. To overcome such issue, the calibration method relies on a second assumption: another linear relation exists between biomarker-derived ( $M_i$ ) and true intake data:

$$M_i = \mu_M + \beta_M X_i + \epsilon_{M,i} \quad (2)$$

Further, it is assumed that biomarker-derived estimates of intake provide a more reliable representation of true intakes, setting  $\beta_M = 1$  and  $\mu_M = 0$ , thus following the classical measurement error model REF here Carroll book 2006. Subject-specific errors  $\epsilon_{M,i}$  are assumed to be independently distributed, according to a Gaussian distribution with zero mean and equal variance  $\sigma_M^2$ .

As the true intakes  $X_i$  cannot be observed, the calibration method combines equations (1) and (2) to derive estimates of such quantities. Specifically, given the assumptions outlined hold,  $X_i$  may be conditionally predicted given biomarker-derived estimates of intakes and self-reported intakes in two steps:

$$\begin{aligned} M_i &= \beta_0 + \beta_1 W_i + \epsilon_i, & i = 1, \dots, n \\ \hat{X}_i &= \hat{\beta}_0 + \hat{\beta}_1 W_i, & i = 1, \dots, n \end{aligned}$$

The errors  $\epsilon_i$  are within-subject random effects, distributed independently according to a standard Gaussian distribution. The above equations are the so-called “calibration equations” and the predicted true intakes  $\widehat{X}_i$  are the calibrated self-reported intakes. Indeed, the estimates of the fixed effect population level parameters  $\widehat{\beta}_0$  and  $\widehat{\beta}_1$  may be used to correct self-reported data and get more reliable estimates of true intake; note that such correction incorporates the information coming from biomarker-derived data, but it does not directly depend on it. (4) provides detail on the construction, implementation and theoretical justification of the calibration method.

### Alternative model formulations

Intake data is often characterized by a large presence of zero values. Depending on whether the measure of intake is self-reported or biomarker-derived, zero values either indicate non-consumers (subjects that reported null consumption of a given food) or null biomarker-derived estimates. There may be a partial mismatch between subjects self-reporting zero intake values and those having zero biomarker-derived estimates of intake values. When subject-specific covariates have been collected, several approaches, eg. (5), have been proposed to incorporate them in the modelling framework and to try modelling “true” and “false” (those misreporting) non-consumers. In the absence of covariate data such zero self-reported intakes are more challenging to address, but modelling frameworks designed for zero-inflated data can be exploited eg. (6). In particular, the calibration equations may be re-defined according to either a Tobit model (7) or a two-part model (8); such models are designed specifically for a regression framework where the response variable,  $M_i$ , is semi-continuous. Then, the calibration equations may be rewritten as:

$$\begin{aligned} M_i &= f_l(W_i) + \epsilon_i, \quad i = 1, \dots, n \\ \widehat{X}_i &= \widehat{f}_l(W_i) \end{aligned}$$

The function  $f_l(\cdot)$  is a link function, corresponding to the linear regression, the Tobit regression or the two-part regression, respectively for  $l = 1, 2, 3$ .

### Biomarker transformation

The assumption of Gaussian errors  $\epsilon_i$  may not be completely satisfied in empirical studies, and some skewness may be found in observed intake data. As can be seen from the summary statistics provided in the paper, this is the case for the observed mean daily self-reported intakes and biomarker-derived estimates of citrus intake. In particular, when deriving calibration equations for citrus intake, the skewness in the distribution of biomarker-derived estimates of intake may imply non-Normality of the errors  $\epsilon_i$ . To try overcome this potential issue, in the context of episodically consumed foods (5) propose to transform the response variable using the Box-Cox family of transformations (9),

$$M_i^* = BC(M_i, \lambda) = \frac{M_i^\lambda - 1}{\lambda}, \quad i = 1, \dots, n$$

so that the distribution of the transformed marker data  $M_i^*$  would more closely resemble that of a Gaussian. Using a Box-Cox transformation of biomarker-derived intakes, the calibration equations may be defined as:

$$\begin{aligned} M_i^* &= f_l(W_i) + \epsilon_i, \quad i = 1, \dots, n \\ \widehat{X}_i^* &= \widehat{f}_l(W_i) \end{aligned}$$

where  $\widehat{X}_i^*$  is the transformed predicted true intake and  $f_l(\cdot)$  is any given model specification for the calibration equations. Predicted intakes on the original scale may be found by computing the inverse of the Box-Cox transformation,  $\widehat{X}_i = BC^{-1}(\widehat{X}_i^*, \lambda)$ , using the Delta method (10):

$$\hat{X}_i = (\lambda \hat{X}_i^* + 1)^{\frac{1}{\lambda}} \left[ 1 + \frac{(1 - \lambda)}{2(\lambda \hat{X}_i^* + 1)^2} \right]$$

Note that when biomarker-derived intake data are transformed using a Box-Cox transformation, the two-part and Tobit model can no longer be estimated, as structural zeros “disappear” in the transformed variable.

A different transformation that attempts to normalize the data while leaving the structural zeros unchanged is the log-transformation  $M_i^* = \log(M_i + 1)$ . This type of log-transformation is widely used in dose-response study data, where there is often the need to rescale the data without losing information on the structural zero values. Indeed, for  $M_i = 0$ ,  $M_i^* = \log(1) = 0$ . As for Box-Cox transformed intake data, when using log-transformed biomarker-derived estimates of intake to estimate calibration equations, the computed calibrated intakes are defined on a new, transformed, scale. Calibrated intake on the original scale may be obtained using the back transformation  $\hat{X}_i = \frac{3}{2} \exp(\hat{X}_i^*) - 1$ .

### Estimating calibration equations: cross validation and mean squared error

In order to select the optimal model specification,  $\hat{f}_l(\cdot)$ ,  $l = 1, 2, 3$ , and data transformation (original scale, Box-Cox or log-transformation), different calibration equation models have been estimated, resulting from all possible combinations of model specifications and data transformations. For the Box-Cox transformation, parameter values  $\lambda = 0.05, 0.1, 0.15, \dots, 1.2$  have been considered, as these yield the most Gaussian-like distributions of transformed biomarker derived estimates of intake. As a result, 30 different models have been estimated. To obtain stable estimates of the calibration equations, each model  $\mathcal{M}_k$ ,  $k = 1, \dots, K$  ( $K = 30$ ), has been estimated using “leave-p-out cross-validation”, with “p” corresponding to 20% of the data. Such operation has been repeated  $T = 5000$  times,  $t = 1, \dots, T$ , each time randomly splitting the observations in  $A_n$  into a training set (containing 80% of the data) and a test set. For each  $\mathcal{M}_k$  model, calibration equations have been estimated using training data only; the corresponding estimated coefficients are then used to compute calibrated intakes from the self-reported intakes in the test set:

$$M_{i,k,t}^{train} = f_l(W_{i,k,t}^{train}) + \epsilon_{i,k,t}^{train} \quad (3)$$

$$\hat{X}_{i,k,t}^{test} = \hat{f}(W_{i,k,t}^{test}) \quad (4)$$

When using transformed biomarker-derived estimates of intake, in equations (3) and (4),  $M_{i,k,t}^{train}$  may be substituted with the transformed-intake  $M_{i,k,t}^{train,*}$ , and  $\hat{X}_{i,k,t}^{test}$  would then be the transformed-calibrated intake  $\hat{X}_{i,k,t}^{test,*}$ .

To compare calibrated intakes obtained with different calibration equations models  $\mathcal{M}_k$ ,  $k = 1, \dots, K$ , mean squared errors are computed, between the known biomarker-derived intake values  $M_{i,k,t}^{test}$  and the calibrated intake values  $\hat{X}_{i,k,t}^{test}$ :

$$MSE_{k,t} = \frac{\sum_{i=1}^{0.20n} (M_{i,k,t}^{test} - \hat{X}_{i,k,t}^{test})^2}{0.20n}$$

Different models may then be compared using average mean square errors  $\overline{MSE}_k = \sum_{t=1}^T MSE_{k,t} / T$ , with lower values corresponding to better models. Thus, the optimal model  $\mathcal{M}_{k^*}$  may be found as  $k^* = \arg \min_k \overline{MSE}_k$ .

## SIMULATION STUDY

The calibration equations framework is built upon the assumption that both self-reported intakes and biomarker-derived estimates of intake are linearly related to the unknown true intakes. Normally, biomarker-derived intakes are collected only for a sub-sample of the population of interest, due to elevated extraction costs. Indeed, in the mean daily citrus intake data the population of interest counted  $N = 1500$  subjects, but biomarker-derived intakes had been collected only for  $n = 565$  of them, approximately one third of the total. With large population studies, it may be infeasible to collect biomarker-derived intakes for a large subset of the subjects. Biomarker-derived data would be then collected only for a small number of  $n$  subjects,  $n \ll N$ . However, deriving calibration equations using  $n \ll N$  observations may result in unstable and potentially unreliable intake predictions. Hence, how large, or small,  $n$  is required to be in order to obtain stable predictions is of key interest. To investigate this issue, a simulation study has been constructed. Assuming the population of interest has  $N = 100,000$  subjects,  $P = 100$  different sets of true, self-reported and biomarker-derived mean daily citrus intakes have been simulated according to the calibration equations framework:

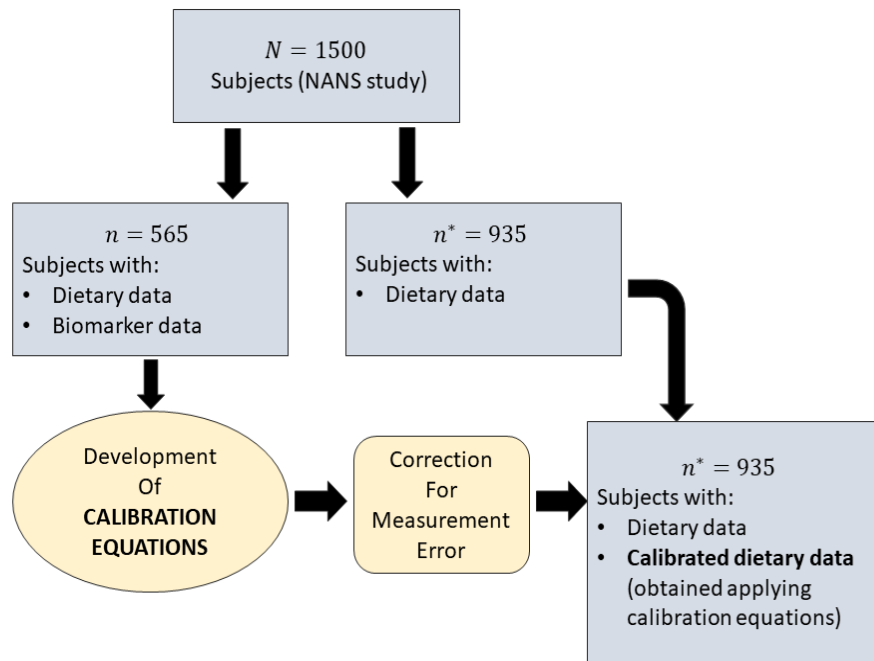
$$\begin{aligned} W_{i,p} &= \mu_W + \beta_W X_{i,p} + \epsilon_{i,p,W} \\ M_{i,p} &= X_{i,p} + \epsilon_{i,p,M} \end{aligned}$$

True intakes  $X_{i,p}$  have been simulated from a Gaussian distribution  $X_{i,p} \sim \mathcal{N}(\mu_X, \sigma_X^2)$ , where  $\mu_X = 81$  and  $\sigma_X^2 = (66)^2$  have been fixed to reflect the estimated mean and variance of the calibrated mean daily citrus intakes. The random errors,  $\epsilon_{i,p,W}$  and  $\epsilon_{i,p,M}$  are assumed to be normally distributed with zero mean and variances  $\sigma_W^2 = \alpha_W \sigma_X^2$  and  $\sigma_M^2 = \alpha_M \sigma_X^2$ . The  $\alpha_W$  and  $\alpha_M$  scale parameters regulate the variation of the self-reported and biomarker-derived intake data around the true intake. These scale parameters have been set to  $\alpha_W, \alpha_M = (0.5, 2)$ , to represent both clean ( $\alpha_W, \alpha_M = 0.5$ ) and noisy ( $\alpha_W, \alpha_M = 2$ ) intake data. The  $\beta_W$  parameter measures the strength of the relation between self-reported and true intake data. To represent different levels of quality for self-reported data, this parameter has been set to  $\beta_W = (0.1, 0.5, 0.8, 1)$ , with  $\beta_W = 0.1$  and  $\beta_W = 1$  denoting poor and good data quality, respectively. Finally, the offset parameter  $\mu_W$  measures the average difference between self-reported and true intake data. Here this parameter has been set to  $\mu_W = 0$ , as similar values were observed for the average self-reported and calibrated mean-daily intakes.

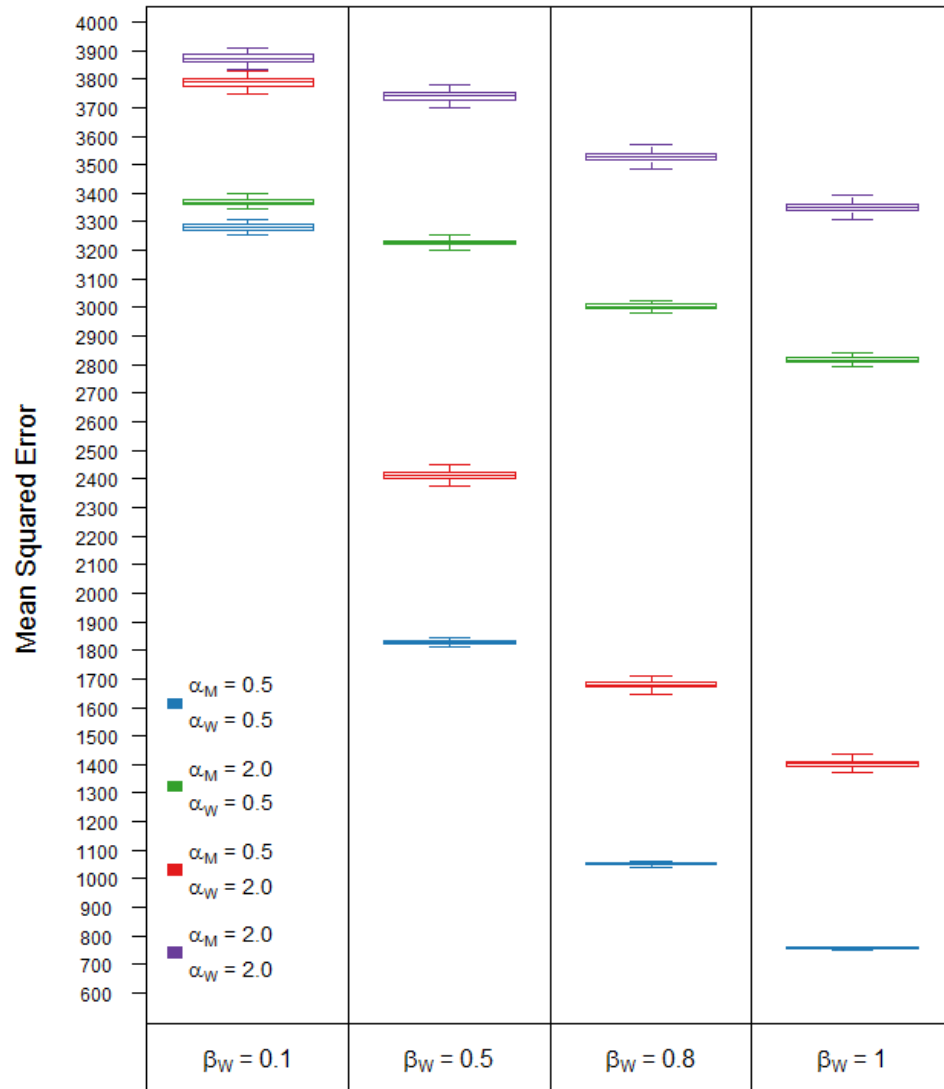
Different combinations of  $\beta_W$ ,  $\alpha_W$  and  $\alpha_M$  parameter values give rise to sixteen different simulation scenarios. For each scenario, it has been further assumed that biomarker-derived estimates of intake have been collected for a subsample of  $n_g$  subjects, with  $n_g = gN$  and  $g = (0.01, 0.02, \dots, 0.99, 1)$ . For each simulation scenario and for each value of  $n_g$  the calibration equations are estimated, using the optimal model from the citrus intake data analysis i.e. the linear model with the original-scale biomarker-derived intakes. As in the citrus intake analysis, calibration equations have been estimated using the “leave-p-out cross-validation” method. The standard deviations of the cross-validation estimates of the calibration equation parameters,  $sd(\hat{\beta}_0)_{\alpha_W, \alpha_M, \beta_W, n_g, p}$  and  $sd(\hat{\beta}_1)_{\alpha_W, \alpha_M, \beta_W, n_g, p}$  are then computed. Finally, these standard deviations are averaged across the  $P$  different simulated datasets giving  $sd(\hat{\beta}_0)_{\alpha_W, \alpha_M, \beta_W, n_g}$  and  $sd(\hat{\beta}_1)_{\alpha_W, \alpha_M, \beta_W, n_g}$ . For each scenario, these averaged standard deviations are explored for increasing number of biomarker-derived intake estimates  $n_g$ . Small decrements of the standard deviations with increasing  $n_g$  would indicate no improvement in parameter estimate stability as a result of, expensively, collecting greater portions of biomarker data.

### Simulation study results

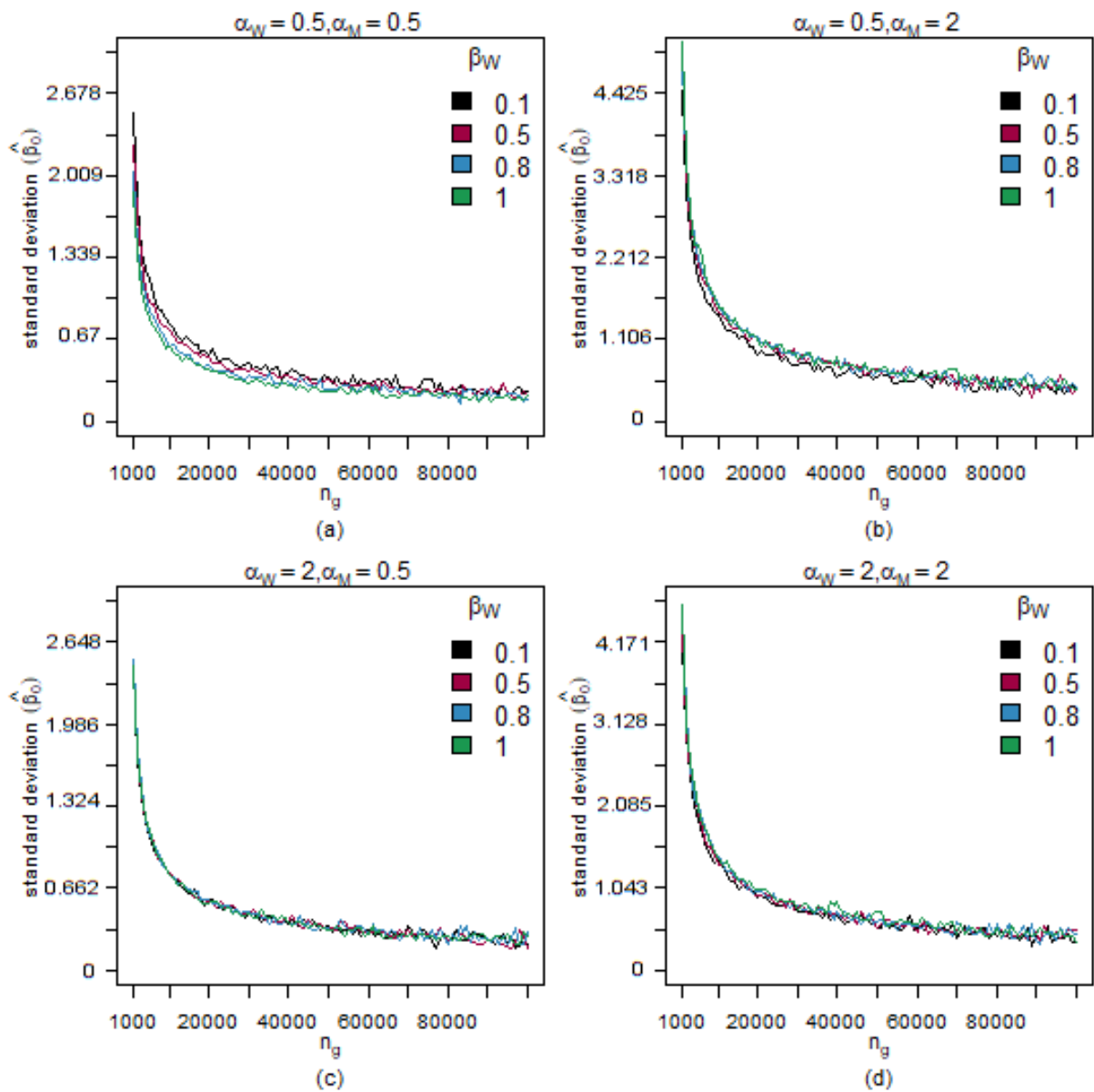
**Supplementary Figure 5** illustrates the distributions of the Mean Squared Errors computed in the different simulation scenarios, between simulated and calibrated mean daily intakes. The MSE increases with decreasing data quality and with increasing variability of the self-reported and biomarker-derived intakes around the true intakes. Notice that when  $\beta_W = 0.1$ , the MSE is not much affected by the noise level  $(\alpha_W, \alpha_M)$ , suggesting that the quality of self-reported data is quite relevant. **Supplementary Figure 6** reports the standard deviations of the estimated intercept  $\beta_0$  in the different scenarios, while **Supplementary Figure 7** those of the estimated slope coefficient  $\beta_1$  in the intermediate scenarios. Standard deviations tend to stabilize quickly, independently of the scenario considered.



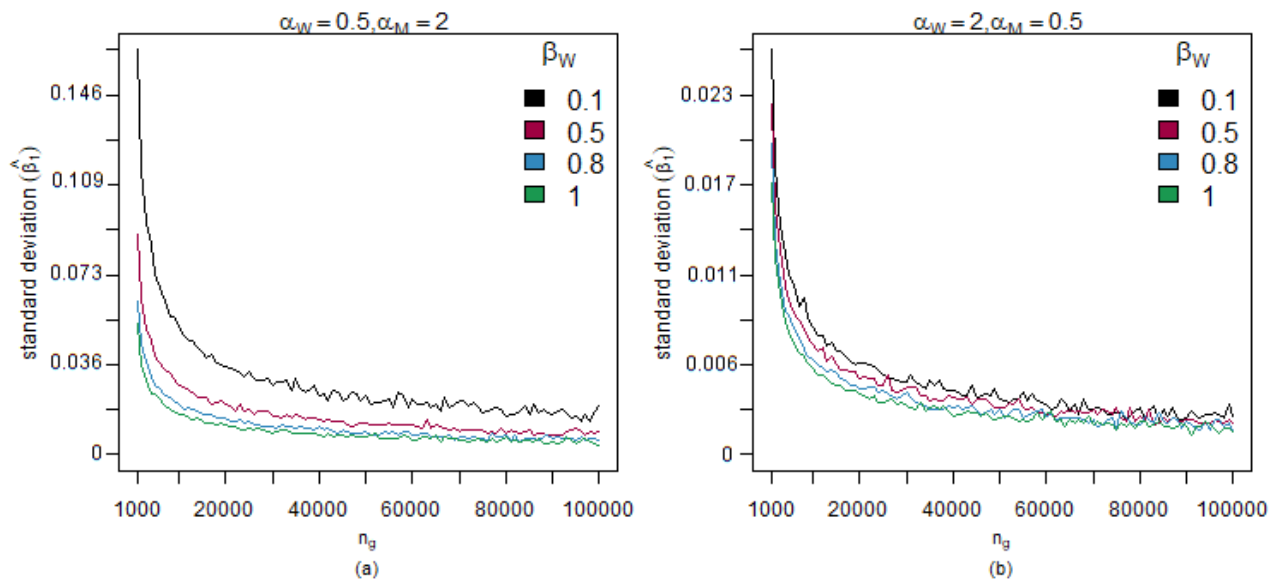
Supplementary Figure 4 Flow chart describing the selection process for the data used to develop calibration equations.



Supplementary Figure 5 Simulations. This figure presents the distributions of the MSE computed in the different simulation scenarios, between simulated and calibrated mean daily intakes. As expected, the error is lower when both self-reported and biomarker-derived intakes have little variability around true intakes, and it decreases with increasing values of  $\beta_W$ . However, when self-reported intake data is of “poor” quality ( $\beta_W = 0.1$ ), the MSE is not much affected by the noise level ( $\alpha_W, \alpha_M$ ).



Supplementary Figure 6 Standard deviations of the estimated intercept  $\beta_0$  in the different scenarios, for different sample sizes  $n_g$  and different values of the  $\beta_W$  parameter. As for the estimated slope coefficient  $\beta_1$ , standard deviations tend to stabilize quickly, independently of the scenario considered.



Supplementary Figure 7 Standard deviations of the estimated slope coefficient  $\beta_1$  in the two intermediate scenarios, for different sample sizes  $n_g$  and different values of the  $\beta_W$  parameter. As for the other scenarios, standard deviations tend to stabilize quickly, even with “poor” self-reported data quality, that is for  $\beta_W = 0.1$ .

## References

1. Prentice RL, Pettinger M, Tinker LF, Huang Y, Thomson CA, Johnson KC, Beasley J, Anderson G, Shikany JM, Chlebowski RT, et al. Regression calibration in nutritional epidemiology: example of fat density and total energy in relationship to postmenopausal breast cancer. *Am J Epidemiol* 2013;178(11):1663-72. doi: 10.1093/aje/kwt198.
2. Prentice RL, Huang Y, Kuller LH, Tinker LF, Horn LV, Stefanick ML, Sarto G, Ockene J, Johnson KC. Biomarker-calibrated energy and protein consumption and cardiovascular disease risk among postmenopausal women. *Epidemiology* 2011;22(2):170-9. doi: 10.1097/EDE.0b013e31820839bc.
3. Huang Y, Van Horn L, Tinker LF, Neuhouser ML, Carbone L, Mossavar-Rahmani Y, Thomas F, Prentice RL. Measurement error corrected sodium and potassium intake estimation using 24-hour urinary excretion. *Hypertension* 2014;63(2):238-44. doi: 10.1161/HYPERTENSIONAHA.113.02218.
4. Gormley IC, Bai Y, Brennan L. Combining biomarker and self-report dietary intake data: a review of the state of the art and an overview of concepts. *Stat Methods Med Res* (in press) 2019.

5. Kipnis V, Midthune D, Buckman DW, Dodd KW, Guenther PM, Krebs-Smith SM, Subar AF, Tooze JA, Carroll RJ, Freedman LS. Modeling data with excess zeros and measurement error: application to evaluating relationships between episodically consumed foods and health outcomes. *Biometrics* 2009;65(4):1003-10. doi: 10.1111/j.1541-0420.2009.01223.x.
6. Min Y, Agresti A. Modeling nonnegative data with clumping at zero: a survey. *JIRSS* 2002;1(1):7-33.
7. Tobin J. Estimation of relationships for limited dependent variables: 1. Introduction. *Econometrica (pre-1986)* 1958;26(1):24.
8. Duan N, Manning WG, Morris CN, Newhouse JP. A comparison of alternative models for the demand for medical care. *J Bus Econ Stat* 1983;1(2):115-26. doi: 10.1080/07350015.1983.10509330.
9. Box GEP, Cox DR. An analysis of transformations. *J R Stat Soc Series B Stat Methodol* 1964;26(2):211-43. doi: doi:10.1111/j.2517-6161.1964.tb00553.x.
10. van der Vaart A.W. The Delta-method. In *Weak Convergence and Empirical Processes*. New York: Springer, 1996.